# Biometrical Modeling of Twin and Family Data
# Using Standard Mixed Model Software

**S. Rabe-Hesketh,[1,*] A. Skrondal,[2] and H. K. Gjessing[3]**

[1]Graduate School of Education & Graduate Group in Biostatistics, University of California,
Berkeley, California 94720, U.S.A. & Institute of Education, University of London, London WC1H OAL, U.K.
[2]Department of Statistics & Methodology Institute, London School of Economics, London, WC2A 2AE, U.K.
& Division of Epidemiology, Norwegian Institute of Public Health, 0403 Oslo, Norway
[3]Division of Epidemiology, Norwegian Institute of Public Health, 0403 Oslo, Norway & Section of Epidemiology
and Medical Statistics, University of Bergen, 5020 Bergen, Norway
[*]*email:* sophiarh@berkeley.edu

SUMMARY. Biometrical genetic modeling of twin or other family data can be used to decompose the variance of an observed response or 'phenotype' into genetic and environmental components. Convenient parameterizations requiring few random effects are proposed, which allow such models to be estimated using widely available software for linear mixed models (continuous phenotypes) or generalized linear mixed models (categorical phenotypes). We illustrate the proposed approach by modeling family data on the continuous phenotype birth weight and twin data on the dichotomous phenotype depression. The example data sets and commands for `Stata` and `R/S-PLUS` are available at the *Biometrics* website.

KEY WORDS: Biometrical genetics; equivalence; family design; generalized linear mixed model; `gllamm`; multilevel model; twin design.

## 1. Introduction

An increasing number of genes have been identified as risk or preventive factors for different phenotypes, that is, the observable characteristics of an individual. However, such knowledge is scarce for many phenotypes and consequently classical biometrical genetics still has an important role to play. In this approach, individual genes are not identified and one instead attempts to decompose the phenotypic variance into genetic and environmental components. The underlying variance components model is specified based on the degree of genetic kinship among individuals in the study population and on various assumptions regarding the degree of shared environment between the same individuals. Estimated variance components may provide some insight into the etiology of diseases but perhaps more importantly direct future genetic research by indicating phenotypes that are largely genetically determined.

Models for twin and other family designs are typically specified as structural equation models and estimated using specialized software, such as `Mx` (Neale et al., 2004). In this multivariate approach, the responses for the family members are treated as different *variables* and the families as the units of analysis. Here we show how the models can be formulated as mixed effects models where the responses for the family members are treated as responses on different *units* nested within families. An important advantage of the mixed models approach is that it is familiar to statisticians and widely available in standard statistical software such as `Stata`, `S-PLUS`, `R`, `SPSS`, and `SAS`, as well as stand-alone software such as `MLwiN` and `HLM`.

One challenge when using the mixed model approach is to impose the restrictions on the covariance matrix of the responses predicated by genetic theory. Pawitan et al. (2004) achieve this using their own custom-made program to estimate nonstandard generalized linear mixed models allowing correlations among random effects to depend on covariates. In contrast, Guo and Wang (2002) ignore the restrictions in order to use standard software for linear mixed models. van den Oord (2001) uses standard software for linear mixed models to impose restrictions, but his approach has several disadvantages: First, complicated parameter constraints are needed. Second, the parameters of the model do not relate simply to the parameters of interest. Third, the models are not proper statistical models because they require theoretically impossible restrictions such as random effects having a zero variance but nonzero covariances with other random effects. Our parameterizations avoid all these problems. In particular, avoiding the third problem means that the models can be estimated using software that forces covariance matrices to be positive semidefinite, as is the case for most software implementing maximum likelihood estimation for categorical responses (e.g., `SAS PROC NLMIXED` and `gllamm` in `Stata`).

An important contribution of our article is that we suggest parameterizations requiring only a few random effects. This is important when the phenotypes are noncontinuous because in this case the likelihood does not have a closed form and

© 2007, The International Biometric Society

the random effects must be integrated out using for instance numerical integration. McArdle and Prescott (2005) recently suggested a similar parameterization as our 'parameterization 1' for twin models. However, their approach requires four random effects whereas ours requires three (see Section 4.1). For the same model, we moreover propose a 'parameterization 2' requiring only two random effects. We also propose parameterizations for more complex family structures.

The plan of the article is as follows. In Section 2, we introduce some basic ideas of biometrical genetic models before considering the special case of twin data in Section 3. Convenient mixed model parameterizations are proposed for twin models in Section 4 and for some other family designs in Section 5. In Section 6, we demonstrate how these parameterizations can be used for categorical phenotypes or qualitative traits. Applications of the suggested approach to continuous as well as dichotomous phenotypes are considered in Section 7. The article concludes with a brief discussion in Section 8.

## 2. Biometrical Genetic Models
In Fisher's (1918) model for polygenic effects the genetic component of a phenotype can be divided into two components. The *additive genetic* component represents the main effects of individual alleles on the phenotype, which are transmissible from parents to offspring. The *dominance genetic* component results from interactions between alleles at single loci, which contribute to the covariance only between relatives who can share a genotype identical by descent (derived from the same parental allele) such as full siblings and twins.

In addition to these effects of 'nature,' 'nurture' obviously also has a role to play. The *common environment* component represents environmental influences shared by siblings reared in the same family. These shared experiences make siblings reared in the same family more alike than siblings reared in different families. In contrast, the *unique environment* component refers to experiences that affect individual siblings (not shared) and make them dissimilar.

We refer to Falconer and MacKay (1996) and Khoury, Beaty, and Cohen (1993, Chapter 7) for motivation of the genetic models. These treatments also discuss the assumptions on which they are based, such as Hardy–Weinberg equilibrium, no epistasis (interactions between alleles at different loci), absence of gene-environment interaction, and random mating. Elston (2001) provides an excellent explanation of important concepts in genetics and a translation between the terminologies of genetics and statistics.

The ACDE model decomposes the total variance of the phenotype into four components, due to Additive genetic, Common environment, Dominance genetic, and unique Environment (ACDE) effects. The model for individual $i$ in family $j$ can be written as an error components model

$$y_{ij} = \mu + A_{ij} + D_{ij} + C_{ij} + \epsilon_{ij}, \tag{1}$$

where $\mu$ is the overall mean, $A_{ij} \sim N(0, \sigma_A^2)$ is an *additive genetic* component, $D_{ij} \sim N(0, \sigma_D^2)$ a *dominance genetic* component, $C_{ij} \sim N(0, \sigma_C^2)$ a *common environment* component, and $\epsilon_{ij} \sim N(0, \sigma_E^2)$ a *unique environment* component.

These four 'error' components are assumed to be mutually independent so that the total variance is the sum of the variance components

$$\text{var}(y_{ij}) = \sigma_A^2 + \sigma_D^2 + \sigma_C^2 + \sigma_E^2.$$

Alternatively, the model can be written as a path model,

$$y_{ij} = \mu + aA_{ij}^* + dD_{ij}^* + cC_{ij}^* + e\epsilon_{ij}^*, \tag{2}$$

where $A_{ij}^*$, $D_{ij}^*$, $C_{ij}^*$, and $\epsilon_{ij}^*$ are mutually independent with standard normal distributions and $a, d, c,$ and $e$ are 'path coefficients.' The squared path coefficients correspond to the variance components in the error components model; $a^2 = \sigma_A^2$, $c^2 = \sigma_C^2$, $d^2 = \sigma_D^2$, and $e^2 = \sigma_E^2$.

## 3. Biometrical Genetic Models for Twin Data
The twin design is a commonly used family design because it is both simple and powerful. Data on the phenotype are obtained for both monozygotic ('identical,' MZ) and dizygotic ('nonidentical,' DZ) twin-pairs. The basic idea is that MZ twins share all their genes and consequently become more similar to each other than DZ twins who have only half their genes identical by descent. Assuming that MZ and DZ twins experience the same degree of similarity in their environments (the equal environment assumption), any excess similarity between MZ twins must be due to the greater proportion of genes shared by MZ twins than by DZ twins.

To represent the covariance structure of the genetic components in a compact manner, we consider two unrelated twin-pairs $j = 1, 2$ with twins $i = 1, 2$ in each pair, where the first pair is MZ and the second DZ. Let the corresponding additive and dominance genetic components be denoted $\mathbf{A} = (A_{11}, A_{21}, A_{12}, A_{22})'$ and $\mathbf{D} = (D_{11}, D_{21}, D_{12}, D_{22})'$, respectively. Analogously, the vectors of common and unique environment components are defined as $\mathbf{C} = (C_{11}, C_{21}, C_{12}, C_{22})'$ and $\mathbf{E} = (E_{11}, E_{21}, E_{12}, E_{22})'$, respectively.

According to genetic theory (e.g., Neale and Maes, 2004), the covariance matrices for the components are

$$\text{Cov}(\mathbf{A}) = \sigma_A^2 \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1/2 \\ 0 & 0 & 1/2 & 1 \end{bmatrix},$$

$$\text{Cov}(\mathbf{D}) = \sigma_D^2 \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1/4 \\ 0 & 0 & 1/4 & 1 \end{bmatrix},$$

$$\text{Cov}(\mathbf{C}) = \sigma_A^2 \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix},$$

$$\text{Cov}(\mathbf{E}) = \sigma_E^2 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \sigma_E^2 \mathbf{I}.$$

Defining the vector of phenotypes for the two twin-pairs as $\mathbf{y} = (y_{11}, y_{21}, y_{12}, y_{22})'$, the covariance structure for the phenotypes under the ACDE model becomes

$$\text{Cov}(\mathbf{y}) = \text{Cov}(\mathbf{A}) + \text{Cov}(\mathbf{D}) + \text{Cov}(\mathbf{C}) + \text{Cov}(\mathbf{E}). \quad (3)$$

Unfortunately, the ACDE model is not identified unless we either have information regarding twins who are reared apart or regarding other relatives. In practice, two special cases of the ACDE model, which are identified from data on MZ and DZ twins are usually considered; the ACE and ADE models where the letters refer to the random effects included in the model. Further special cases include the AE model that is nested in both the ACE and ADE models and the CE model that is nested in the ACE model.

## 4. Parameterization of Twin Models as Linear Mixed Models

Although the twin models look like ordinary random effects models, an unusual feature is that some of the covariance matrices of the random effects depend on the type of twin, a covariate at the twin-pair level. However, most software for linear mixed models does not allow covariance matrices of random effects to be specified as functions of covariates. Furthermore, most software does not permit nonlinear constraints on model parameters. With these limitations in mind, we propose two parameterizations that allow the twin models to be estimated in standard software.

### 4.1 *Parameterization 1*

In the first parameterizations, the same basic trick is used for ACE and ADE models.

**ACE model:** The ACE model can be parameterized with three random effects as follows:

$$y_{ij} = \mu + \left\{ a_{ij}^{(2)} \left[ \sqrt{\tfrac{1}{2}} \overline{M}_j \right] + a_j^{(3)} \left[ M_j + \sqrt{\tfrac{1}{2}} \overline{M}_j \right] \right\} + c_j^{(3)} + \epsilon_{ij}, \quad (4)$$

where $M_j$ is a dummy variable for MZ twins and $\overline{M}_j = 1 - M_j$ a dummy variable for DZ twins. The setup is shown for a hypothetical data set under 'Param. 1' in Table 1, where the first twin-pair is MZ whereas the second is DZ.

$a_{ij}^{(2)}$ is a random coefficient at the individual level whereas $a_j^{(3)}$ and $c_j^{(3)}$ are random coefficients at the twin-pair level. We have used superscripts to denote the levels at which random terms vary apart from $\epsilon_{ij}$, which always varies at level 1. The lowest-level random effects are usually considered to be at level 2, here the individual twins $i$, so that twin-pairs $j$ are at level 3. Here level 2 is nested in level 3 and there is no correspondence between twin 1 in twin-pair 1 and twin

1 in twin-pair 2. For software requiring unique identifiers for different units at a given level (even when they belong to different higher-level units), $i$ could be consecutive integers or the values $i'$ and $j'$ in the table can be used.

All components are specified as mutually uncorrelated and the components $a_{ij}^{(2)}$ and $a_j^{(3)}$ have equal variances

$$\text{Var}\big(a_{ij}^{(2)}\big) = \text{Var}\big(a_j^{(3)}\big) = \sigma_A^2,$$

whereas $\text{Var}(c_j) = \sigma_C^2$ and $\text{Var}(\epsilon_{ij}) = \sigma_E^2$.

The terms in square brackets in (4) are simply linear combinations of dummy variables that can be computed as new variables and specified as having random coefficients. Here the trick is that the additive genetic component for MZ twins is just the shared component $a_j^{(3)}$, producing a variance of $\sigma_A^2$ and a correlation of 1. In contrast, the additive genetic component for DZ twins is a scaled sum $\sqrt{\tfrac{1}{2}}(a_{ij}^{(2)} + a_j^{(3)})$ of the unique component $a_{ij}^{(2)}$ and the shared component $a_j^{(3)}$, producing a total variance of $\sigma_A^2$ and a correlation of $1/2$. This idea of splitting the additive genetic component for DZ twins into a unique or 'within' component and a shared or 'between' component is consistent with the variance components approach of Jinks and Fulker (1970; Table 3).

To confirm that the parameterization is correct, we recommend that the model-implied covariance structure be compared with the one prescribed by genetic theory. Letting $\mathbf{1}$ denote a unit column vector and $\mathbf{I}$ an identity matrix, the model can be written in matrix form as

$$\mathbf{y} = \mathbf{1}\mu + \underbrace{\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{\tfrac{1}{2}} & 0 \\ 0 & 0 & 0 & \sqrt{\tfrac{1}{2}} \end{bmatrix}}_{\mathbf{Z}_1} \begin{bmatrix} a_{11}^{(2)} \\ a_{21}^{(2)} \\ a_{12}^{(2)} \\ a_{22}^{(2)} \end{bmatrix} + \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & \sqrt{\tfrac{1}{2}} \\ 0 & \sqrt{\tfrac{1}{2}} \end{bmatrix}}_{\mathbf{Z}_2} \begin{bmatrix} a_1^{(3)} \\ a_2^{(3)} \end{bmatrix}$$

$$+ \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}}_{\mathbf{Z}_3} \begin{bmatrix} c_1^{(3)} \\ c_2^{(3)} \end{bmatrix} + \mathbf{I} \begin{bmatrix} \epsilon_{11} \\ \epsilon_{21} \\ \epsilon_{12} \\ \epsilon_{22} \end{bmatrix},$$

and the covariance structure becomes

$$\text{Cov}(\mathbf{y}) = \sigma_A^2(\mathbf{Z}_1\mathbf{Z}_1' + \mathbf{Z}_2\mathbf{Z}_2') + \sigma_C^2\mathbf{Z}_3\mathbf{Z}_3' + \sigma_E^2\mathbf{I}$$

as required.

**Table 1**
*Hypothetical data setup for the two suggested parameterizations for twin models*

| Twin-pair $j$ | Twin $i$ | Phenotype $y_{ij}$ | Param. 1 | | Param. 2 $k$ | Alt. indices | | |
|---|---|---|---|---|---|---|---|---|
| | | | $M_j$ | $\overline{M}_j = 1 - M_j$ | | $j'$ | $i'$ | $k'$ |
| 1 | 1 | $y_{11}$ | 1 | 0 | 1 | 3 | 1 | 3 |
| 1 | 2 | $y_{21}$ | 1 | 0 | 1 | 3 | 2 | 3 |
| 2 | 1 | $y_{12}$ | 0 | 1 | 1 | 6 | 4 | 4 |
| 2 | 2 | $y_{22}$ | 0 | 1 | 2 | 6 | 5 | 5 |

The AE model is simply obtained by omitting the penultimate term and the CE model by omitting the second and third terms from the above model.

McArdle and Prescott (2005) propose a similar parameterization that requires two random effects in place of our single random effect $a_{ij}^{(2)}$. Defining dummy variables $T_{1i}$ for $i = 1$ and and $T_{2i}$ for $i = 2$, their model can be written as

$$y_{ij} = \mu + \left\{ a_{1ij}^{(2)} \left[ \sqrt{\frac{1}{2}} \overline{M}_j T_{1i} \right] + a_{2ij}^{(2)} \left[ \sqrt{\frac{1}{2}} \overline{M}_j T_{2i} \right] \right. $$
$$\left. + a_j^{(3)} \left[ M_j + \sqrt{\frac{1}{2}} \overline{M}_j \right] \right\} + c_j^{(3)} + \epsilon_{ij},$$

where $\operatorname{var}(a_{1ij}^{(2)}) = \operatorname{var}(a_{2ij}^{(2)}) = \operatorname{var}(a_j^{(3)}) = \sigma_A^2$. Due to the variance constraint, and because the random effects at level 2 take on different values for different $i$ ($i = 1$, 2) the above model is equivalent to our parameterization 1 but computationally more demanding.

**ADE model:** Our suggested parameterization for the ADE model requires four random effects:

$$y_{ij} = \mu + \left\{ a_{ij}^{(2)} \left[ \sqrt{\frac{1}{2}} \overline{M}_j \right] + a_j^{(3)} \left[ M_j + \sqrt{\frac{1}{2}} \overline{M}_j \right] \right\}$$
$$ + \left\{ d_{ij}^{(2)} \left[ \sqrt{\frac{3}{4}} \overline{M}_j \right] + d_j^{(3)} \left[ M_j + \sqrt{\frac{1}{4}} \overline{M}_j \right] \right\} + \epsilon_{ij}, \quad (5)$$

where all components are specified as mutually uncorrelated and

$$\operatorname{Var}\left(a_{ij}^{(2)}\right) = \operatorname{Var}\left(a_j^{(3)}\right) = \sigma_A^2, \quad \operatorname{Var}\left(d_{ij}^{(2)}\right) = \operatorname{Var}\left(d_j^{(3)}\right) = \sigma_D^2, \quad \text{and}$$
$$\operatorname{Var}(\epsilon_{ij}) = \sigma_E^2.$$

For two twin-pairs, the first MZ and the second DZ, the ADE model can be written as

$$\mathbf{y} = \mathbf{1}\mu + \mathbf{Z}_1 \begin{bmatrix} a_{11}^{(2)} \\ a_{21}^{(2)} \\ a_{12}^{(2)} \\ a_{22}^{(2)} \end{bmatrix} + \mathbf{Z}_2 \begin{bmatrix} a_1^{(3)} \\ a_2^{(3)} \end{bmatrix} + \underbrace{\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{\frac{3}{4}} & 0 \\ 0 & 0 & 0 & \sqrt{\frac{3}{4}} \end{bmatrix}}_{\mathbf{Z}_4} \begin{bmatrix} d_{11}^{(2)} \\ d_{21}^{(2)} \\ d_{12}^{(2)} \\ d_{22}^{(2)} \end{bmatrix}$$

$$ + \underbrace{\begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & \sqrt{\frac{1}{4}} \\ 0 & \sqrt{\frac{1}{4}} \end{bmatrix}}_{\mathbf{Z}_5} \begin{bmatrix} d_1^{(3)} \\ d_2^{(3)} \end{bmatrix} + \mathbf{I} \begin{bmatrix} \epsilon_{11} \\ \epsilon_{21} \\ \epsilon_{12} \\ \epsilon_{22} \end{bmatrix},$$

and the covariance structure becomes

$$\operatorname{Cov}(\mathbf{y}) = \sigma_A^2 (\mathbf{Z}_1 \mathbf{Z}_1' + \mathbf{Z}_2 \mathbf{Z}_2') + \sigma_D^2 (\mathbf{Z}_4 \mathbf{Z}_4' + \mathbf{Z}_5 \mathbf{Z}_5') + \sigma_E^2 \mathbf{I}.$$

A very convenient feature of these parameterizations is the direct correspondence between the required variance components and the model parameters. As a consequence, any software that does not permit negative variance components enforces the useful restriction that the genetic and environmental variance components must all be nonnegative.

## 4.2 *Parameterization 2*

The second parameterization can be used either for ACE or ADE models. We use the generic notation $u$ for random effects when they do not contribute to a single source $(A, C, D, E)$ of variation.

First, define a new cluster identifier $k$ for level 2 that is equal to $j$ (the twin-pair identifier) for MZ twins and $i$ (the twin identifier) for DZ twins. The random effect $u_{kj}^{(2)}$ is, therefore, shared by individuals in the same twin-pair for MZ twins but unique to the individuals for DZ twins. The setup is shown under 'Param. 2' in Table 1. For software not allowing different level-2 units to have the same values of the level-2 identifier even if they belong to different level-3 units, the alternative indices $i'$, $k'$, and $j'$ in Table 1 can be used.

The model can then be parameterized as a three-level random intercept model

$$y_{ikj} = \mu + u_{kj}^{(2)} + u_j^{(3)} + \epsilon_{ikj}, \quad (6)$$

where the error components are mutually uncorrelated. Here the trick is to let the variance of $u_{kj}^{(2)}$ be shared for MZ twins and unique for DZ twins by defining the artificial level 2.

For the ACE model, $\operatorname{Var}(u_{kj}^{(2)}) = \sigma_A^2/2$, $\operatorname{Var}(u_j^{(3)}) = \sigma_C^2 + \sigma_A^2/2$, and $\operatorname{Var}(\epsilon_{ikj}) = \sigma_E^2$. For the ADE model, $\operatorname{Var}(u_{kj}^{(2)}) = \sigma_A^2/2 + 3\sigma_D^2/4$, $\operatorname{Var}(u_j^{(3)}) = \sigma_A^2/2 + \sigma_D^2/4$, and $\operatorname{Var}(\epsilon_{ikj}) = \sigma_E^2$. The AE model is obtained by imposing the restriction $\operatorname{Var}(u_{kj}^{(2)}) = \operatorname{Var}(u_j^{(3)})$ whereas the CE model can be obtained by omitting $u_{kj}^{(2)}$ from (6).

For two unrelated twin-pairs, the first MZ and the second DZ, the model can be written as

$$\mathbf{y} = \mathbf{1}\mu + \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{Z}_6} \begin{bmatrix} u_{11}^{(2)} \\ u_{12}^{(2)} \\ u_{22}^{(2)} \end{bmatrix} + \mathbf{Z}_3 \begin{bmatrix} u_1^{(3)} \\ u_2^{(3)} \end{bmatrix} + \mathbf{I} \begin{bmatrix} \epsilon_{11} \\ \epsilon_{21} \\ \epsilon_{12} \\ \epsilon_{22} \end{bmatrix}.$$

For the ACE model, the covariance structure becomes

$$\operatorname{Cov}(\mathbf{y}) = \left(\sigma_A^2/2\right)(\mathbf{Z}_6 \mathbf{Z}_6' + \mathbf{Z}_3 \mathbf{Z}_3') + \sigma_C^2 \mathbf{Z}_3 \mathbf{Z}_3' + \sigma_E^2 \mathbf{I},$$

and analogously for the ADE model.

An advantage of this parameterization is that it is very simple and does not require any constraints. More importantly, it only requires two random effects thus increasing computational efficiency when the responses are categorical. A disadvantage is that the genetic and environmental variance components of interest are linear combinations of the variances of the model. Another disadvantage of the parameterization is that the model does not impose the restrictions that the genetic and environmental variance components are all nonnegative. However, this can also be seen as an advantage. In the ACE model, a negative estimate of $\sigma_C^2$ is an indication that the model is not appropriate. In parameterizations enforcing nonnegative variance components, the estimate would be zero. In this case, it is common practice to consider the ADE model, which can be obtained from parameterization 2 without the need to re-estimate the model.

## 5. Parameterizations for Some Other Family Designs

### 5.1 *Siblings Plus Cousins*

To represent the covariance structure compactly, we consider two pairs of siblings, all sharing the same grandparents. The parameterization given below also applies to larger families.

The covariance matrices of the additive genetic and common environmental effects become, respectively,

$$\text{Cov}(\mathbf{A}) = \sigma_A^2 \begin{bmatrix} 1 & 1/2 & 1/8 & 1/8 \\ 1/2 & 1 & 1/8 & 1/8 \\ 1/8 & 1/8 & 1 & 1/2 \\ 1/8 & 1/8 & 1/2 & 1 \end{bmatrix} \quad \text{and}$$

$$\text{Cov}(\mathbf{C}) = \sigma_C^2 \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix},$$

see for instance Khoury et al. (1993).

Treating the wider family $k$ as level 3 and the sibling-pairs $j$ as level 2, the ACE model can be parameterized using two random effects as

$$y_{ijk} = \mu + a_k^{(3)} + u_{jk}^{(2)} + \epsilon_{ijk},$$

where $\text{Var}(a_k^{(3)}) = \sigma_A^2/8$, $\text{Var}(u_{jk}^{(2)}) = 3\sigma_A^2/8 + \sigma_C^2$, and $\text{Var}(\epsilon_{ijk}) = \sigma_A^2/2 + \sigma_E^2$. Here the level-3 component $a_k^{(3)}$ produces the required covariances among cousins, whereas the level-2 component $u_{jk}^{(2)}$ produces the additional covariance among siblings due to both the common environment and the closer kinship. The level-1 residual $\epsilon_{ijk}$ represents additional additive genetic variance and the unique environment.

In matrix notation,

$$\mathbf{y}_k = \mathbf{1}\mu + \mathbf{1}a_k^{(3)} + \mathbf{Z}_3\mathbf{u}_k^{(2)} + \mathbf{I}\epsilon_k,$$

where $\mathbf{u}_k^{(2)} = (u_{1k}^{(2)}, u_{2k}^{(2)})'$ and $\epsilon_k = (\epsilon_{11k}, \epsilon_{21k}, \epsilon_{12k}, \epsilon_{22k})'$, so that the covariance structure becomes

$$\text{Cov}(\mathbf{y}_k) = \left(\sigma_A^2/8\right)(\mathbf{1}\mathbf{1}' + 3\mathbf{Z}_3\mathbf{Z}_3' + 4\mathbf{I}) + \sigma_C^2\mathbf{Z}_3\mathbf{Z}_3' + \sigma_E^2\mathbf{I},$$

as required.

### 5.2 *Nuclear Family*

Without loss of generality, we consider a nuclear family with two children. The covariance matrices of the additive genetic and common environmental effects (for mother, father, child1, child2) are, respectively,

$$\text{Cov}(\mathbf{A}) = \sigma_A^2 \begin{bmatrix} 1 & 0 & 1/2 & 1/2 \\ 0 & 1 & 1/2 & 1/2 \\ 1/2 & 1/2 & 1 & 1/2 \\ 1/2 & 1/2 & 1/2 & 1 \end{bmatrix} \quad \text{and}$$

$$\text{Cov}(\mathbf{C}) = \sigma_C^2 \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix},$$

see Khoury et al. (1993).

We use three nested grouping identifiers; $i$ for individuals treated as level 2 (because we define a random effect for individuals), $j$ taking on unique values for each sibling-pair and each parent (an artificial level 3), and $k$ for families at level 4. The model can then be parameterized using four random effects as

$$y_{ijk} = \mu + a_{1k}^{(4)}[M_i + K_i/2] + a_{2k}^{(4)}[F_i + K_i/2]$$
$$+ a_{ijk}^{(2)}[K_i/\sqrt{2}] + c_{jk}^{(3)} + \epsilon_{ijk}, \tag{7}$$

where $M_i$ is a dummy variable for mothers, $F_i$ for fathers, and $K_i$ for children. The first three terms represent the additive genetic component and the corresponding random effects have variance $\sigma_A^2$, whereas $\text{Var}(c_{jk}^{(3)}) = \sigma_C^2$ and $\text{Var}(\epsilon_{ijk}) = \sigma_E^2$. Here $a_{1k}^{(4)}$ and $a_{2k}^{(4)}$ induce the required additive genetic covariances between each parent and each child and among the children. However, the induced variances for the children are only $\sigma_A^2/2$ and the remaining variance $\sigma_A^2/2$ is provided by $a_{ijk}^{(2)}$. The common environmental component $c_{jk}^{(3)}$ is shared among the children but unique for each parent as required due to the artificial level-3 identifier $j$. Finally, the unique environmental component is represented by $\epsilon_{ijk}$.

In matrix notation,

$$\mathbf{y}_k = \mathbf{1}\mu + \underbrace{\begin{bmatrix} 1 \\ 0 \\ 1/2 \\ 1/2 \end{bmatrix}}_{\mathbf{Z}_7} a_{1k}^{(4)} + \underbrace{\begin{bmatrix} 0 \\ 1 \\ 1/2 \\ 1/2 \end{bmatrix}}_{\mathbf{Z}_8} a_{2k}^{(4)}$$

$$+ \underbrace{\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{\tfrac{1}{2}} & 0 \\ 0 & 0 & 0 & \sqrt{\tfrac{1}{2}} \end{bmatrix}}_{\mathbf{Z}_1} \begin{bmatrix} a_{11k}^{(2)} \\ a_{22k}^{(2)} \\ a_{33k}^{(2)} \\ a_{43k}^{(2)} \end{bmatrix} + \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{Z}_9} \begin{bmatrix} c_{1k}^{(3)} \\ c_{2k}^{(3)} \\ c_{3k}^{(3)} \end{bmatrix} + \mathbf{I}\epsilon_k,$$

and

$$\text{Cov}(\mathbf{y}_k) = \sigma_A^2(\mathbf{Z}_7\mathbf{Z}_7' + \mathbf{Z}_8\mathbf{Z}_8' + \mathbf{Z}_1\mathbf{Z}_1') + \sigma_C^2\mathbf{Z}_9\mathbf{Z}_9' + \sigma_E^2\mathbf{I}.$$

For software permitting negative variance component estimates, a parametrization requiring only three random effects is

$$y_{ijk} = \mu + a_{1k}^{(4)}[M_i - K_i/\sqrt{2}] + a_{2k}^{(4)}[F_i - K_i/\sqrt{2}] + u_{jk}^{(3)} + \epsilon_{ij},$$

with $\text{Var}(a_{1k}^{(4)}) = \text{Var}(a_{2k}^{(4)}) = -\sigma_A^2/2$, $\text{Var}(u_{jk}^{(3)}) = \frac{1+\sqrt{2}}{2}\sigma_A^2$, and $\text{Var}(\epsilon_{ijk}) = \sigma_E^2 + \frac{1}{2}\sigma_A^2$.

## 6. Extension to Categorical Phenotypes

For categorical (dichotomous and ordinal) phenotypes, we can assume that the previously considered models hold for an underlying continuous response or *liability* $y_{ij}^*$ (e.g., Falconer and MacKay, 1996) instead of an observed phenotype $y_{ij}$, producing *probit* models.

The phenotype is determined by the liability via a threshold model

$$y_{ij} = s \quad \text{if} \quad \kappa_s < y_{ij}^* \le \kappa_{s+1}, \quad s = 0, \dots, S-1,$$

where $\kappa_s$, $s = 1, \ldots, S - 1$, are fixed unknown thresholds and $\kappa_0 = -\infty$ and $\kappa_S = \infty$. Note that dichotomous phenotypes are simply obtained as the special case where $S = 2$. Similar models can also be used for censored responses such as survival times or responses with ceiling or floor effects. In this case $y_i^*$ is observed unless it is lower than some threshold (left censoring) or higher than some threshold (right censoring).

To illustrate, we consider the ACE twin model for $y_{ij}^*$,

$$y_{ij}^* = A_{ij} + C_j + \epsilon_{ij}, \qquad \text{var}(\epsilon_{ij}) = 1, \tag{8}$$

with $C_j$ replaced by $D_{ij}$ for the ADE model. The mean $\mu$ has been omitted because the threshold model already includes $S - 1$ parameters for the $S$ probabilities.

Given the random effects, the cumulative probability that $y_{ij} \geq s$ in the ACE model is

$$\Pr(y_{ij} \geq s \,|\, A_{ij}, C_j) = \Phi(A_{ij} + C_j - \kappa_s),$$

and similarly for the ADE model.

Using the parameterizations 1 or 2 in (4) or (6), respectively, we obtain three-level generalized linear mixed models (e.g., Rabe-Hesketh and Skrondal, 2007). For parameterization 2 with variance parameters $\psi^{(2)} \equiv \text{Var}(u_{kj}^{(2)})$ and $\psi^{(3)} \equiv \text{Var}(u_j^{(3)})$, the marginal likelihood can be expressed as

$$\begin{aligned}
&\mathrm{L}\big(\psi^{(2)}, \psi^{(3)}, \kappa_1, \ldots, \kappa_{S-1}\big) \\
&= \prod_j \int_{u_j^{(3)}} g\big(u_j^{(3)}; \psi^{(3)}\big) \\
&\quad \times \left[ \prod_k \int_{u_{kj}^{(2)}} g\big(u_{kj}^{(2)}; \psi^{(2)}\big) \right. \\
&\quad \left. \times \left\{ \prod_i f\big(y_{ij} \,|\, u_{kj}^{(2)}, u_j^{(3)}; \kappa_1, \ldots, \kappa_{S-1}\big) \right\} \mathrm{d}u_{kj}^{(2)} \right] \mathrm{d}u_j^{(3)},
\end{aligned}$$

where $f(y_{ij} \,|\, u_{kj}^{(2)}, u_j^{(3)}; \kappa_1, \ldots, \kappa_{S-1})$ is the conditional probability of the observed response $y_{ij}$ given the random effects and $g(\cdot; \psi)$ is a normal density with zero mean and variance $\psi$. The integrals do not have closed forms but can be evaluated by numerical integration (e.g., Rabe-Hesketh, Skrondal, and Pickles, 2005) or Monte Carlo integration.

Some models require three or more random effects, making estimation methods relying on numerical (or Monte Carlo) integration computationally expensive. The most common alternative approach is penalized quasilikelihood (e.g., Breslow and Clayton, 1993) implemented in software such as R, S-PLUS, SAS PROC GLIMMIX, MLwiN, and HLM. However, this approach is likely to perform poorly for family data because it is known to produce biased estimates when cluster sizes are small and/or there is large intra-family dependence (e.g. Rodriguez and Goldman, 1995, 2001; Breslow, 2005). Software implementing maximum likelihood via numerical integration for three-level generalized linear mixed models includes gllamm in Stata (Rabe-Hesketh and Skrondal, 2005) and the stand-alone program aML (Lillard and Panis, 2000). Programs for two-level models such as SAS PROC NLMIXED or MIXOR (Hedeker and Gibbons, 1996) can also sometimes be used, for instance for nuclear family data with a single child per family (see Section 7.1). Currently, only gllamm and SAS

PROC NLMIXED use adaptive quadrature, which is superior to ordinary quadrature. Markov chain Monte Carlo methods are implemented in MLwiN (Browne, Rasbash, and Ng, 2005) and Bugs (Spiegelhalter et al., 1996).

## 7. Examples

### 7.1 *Birthweight: ACE Model with Covariates for Continuous Data on Nuclear Families*

We analyze a random subset of the birth weight data from the Medical Birth Registry of Norway described in Magnus et al. (2001). There are 1000 nuclear families each comprising mother, father, and a single child (not necessarily the only child in the family).

With only a single child per family, the common environment variance is not identified. Further, the random effect $a_{ijk}^{(2)}$ in (7) could be treated either as an individual-specific random effect or as a family-specific random effect because it is multiplied by the dummy $K_i$ for child, which equals '1' for only one individual per family. The model becomes a two-level model for family members $i$ nested in families $j$ with three uncorrelated random coefficients having the same variance,

$$\begin{aligned}
y_{ij} = \mathbf{x}_{ij}'\boldsymbol{\beta} &+ a_{1j}^{(2)}[M_i + K_i/2] + a_{2j}^{(2)}[F_i + K_i/2] \\
&+ a_{3j}^{(2)}[K_i/\sqrt{2}] + \epsilon_{ij},
\end{aligned}$$

where $\mathbf{x}_{ij}$ is a vector of covariates with regression coefficients $\boldsymbol{\beta}$ and the following covariates are included:

(i) [Male]: A dummy variable for being male ($x_{1ij}$)
(ii) [First]: A dummy variable for being the first child ($x_{2ij}$)
(iii) [Midage]: A dummy variable for mother aged 20–35 at time of birth ($x_{3ij}$)
(iv) [Highage]: A dummy variable for mother older than 35 at time of birth ($x_{4ij}$) and
(v) [Birthyr]: Year of birth minus 1967 (earliest birth year in birth registry) ($x_{5ij}$)

The estimates for this model using Stata's xtmixed command or R/S-PLUS's lme function are given in Table 2.

As expected, males weigh more than females, first borns weigh less than subsequent children, birthweight increases with mothers' age, and there is a positive period effect. The additive genetic component has a similar standard deviation

**Table 2**
*Maximum likelihood estimates (in grams) for nuclear family birthweight data*

|  | Estimate | (SE) |
|---|---|---|
| Fixed part | | |
| $\beta_0$ [Constant] | 3461.46 | (34.78) |
| $\beta_1$ [Male] | 158.45 | (17.35) |
| $\beta_2$ [First] | −139.40 | (18.74) |
| $\beta_3$ [Midage] | 57.06 | (31.90) |
| $\beta_4$ [Highage] | 118.86 | (54.67) |
| $\beta_5$ [Birthyr] | 3.63 | (0.69) |
| Random part | | |
| $\sigma_A$ | 315.06 | (16.12) |
| $\sigma_E$ | 365.46 | (12.41) |
| Log likelihood | −22,746.23 | |

as the unique environment component (which represents the sum of common and unique environments).

The 'heritability' $h^2$, defined as the proportion of phenotypic variance explained by the additive genetic factor, is estimated as

$$\widehat{h^2} = \frac{\hat{\sigma}_A^2}{\hat{\sigma}_A^2 + \hat{\sigma}_E^2} = 0.43.$$

To obtain a confidence interval for heritability, we take the inverse standard normal cumulative distributive function (or probit transformation) of the above and compute the corresponding standard error using the delta method. We then form the confidence interval on the probit scale assuming normality and back-transform the confidence limits using the standard normal cumulative distributive function. The resulting 95% confidence interval is from 0.35 to 0.50.

### 7.2 *Depression: ACE and ADE Models for Dichotomous Twin Data*

We analyze data on psychiatric disorders in Caucasian female twin-pairs sampled from the Virginia Twin Registry. The ages of the participants at the time of the interview ranged from 17 to 55. Lifetime psychiatric illness was diagnosed using an adapted version of the Structured Clinical Interview for DSM-III-R Diagnosis. Each member of a twin-pair was interviewed by a different interviewer. The data have been analyzed by Neale and Maes (2004, p. 133).

The parameters were estimated by maximum likelihood with adaptive quadrature using `gllamm` (downloadable from `http://www.gllamm.org`) and are given in Table 3. The estimated common environment variance derived from parameterization 2 for the ACE model is negative. In parameterization 1 the variance is constrained to be nonnegative and the estimate is exactly zero giving the AE model. This model might be selected because it has nearly the same log likelihood as the competing models but with two instead of three parameters.

The observed and predicted numbers of twin-pairs with 0, 1, and 2 cases of depression are given in Table 4 for the AE and ADE models.

For the AE model, the heritability is estimated as

$$\widehat{h^2} = \frac{\hat{\sigma}_A^2}{\hat{\sigma}_A^2 + 1} = 0.43,$$

with approximate 95% confidence interval (via the probit transformation and delta method as before) from 0.31 to 0.53.

**Table 3**
*Maximum likelihood estimates for depression data*

| Parameter | ACE | | ADE | |
|---|---|---|---|---|
| | Param. 1 | Param. 2 | Param. 1 | Param. 2 |
| $\sigma_A^2$ | 0.76 | 0.91 | 0.54 | 0.54 |
| $\sigma_C^2$ | 0.00 | −0.12 | | |
| $\sigma_D^2$ | | | 0.25 | 0.25 |
| $\sigma_E^2$ | 1 | 1 | 1 | 1 |
| Log likelihood | −1257.5 | −1257.4 | −1257.4 | −1257.4 |

**Table 4**
*Observed and model-implied numbers of twin-pairs with* 0, 1, *and* 2 *cases of depression*

| Cases | Observed | Expected AE | Expected ADE |
|---|---|---|---|
| MZ twins | | | |
| 0 | 329 | 311.3 | 312.0 |
| 1 | 178 | 185.5 | 183.9 |
| 2 | 83 | 93.1 | 94.1 |
| DZ twins | | | |
| 0 | 201 | 218.7 | 216.8 |
| 1 | 176 | 165.2 | 168.8 |
| 2 | 63 | 56.0 | 54.4 |

### 8. Discussion

Standard software for linear mixed models (continuous phenotypes) or generalized linear mixed models (categorical phenotypes) can easily be adapted to estimate statistical models for the most common twin and family designs.

An important merit of the mixed models approach is that it is familiar to statisticians and widely available in standard statistical software. Another convenient feature of the mixed model approach is that different family sizes are automatically accommodated as well as responses missing at random. Compared with standard structural equation models, it is also more straightforward to include covariates. Furthermore, when family members are exchangeable, as in twin data, it is not necessary to arbitrarily assign family members to variables and then take extra steps to ensure that the models and estimates are invariant to arbitrary re-assignments. Finally, in models for dichotomous or ordinal phenotypes where maximum likelihood estimation usually requires numerical integration, the dimensionality of the integrals is often much smaller for mixed models than for structural equation models if the families are large.

The structural equation modeling approach also has important merits. First, the required covariance matrices are specified directly rather than indirectly as in mixed models where the covariances are induced by random effects. Second, the biometrical models are often represented as path diagrams (e.g., Neale and Maes, 2004) that can aid understanding. Third, it may be easier to extend the models discussed here to multivariate models for several response variables. A powerful program for both standard and non-standard structural equation modeling called `Mx` can handle missing data, covariates and multilevel models, and has excellent documentation with a focus on biometrical genetics (Neale et al., 2004). This software also allows, for instance, estimation of sibling interaction models, homogamy models and more complex models of parent–child resemblance. These models cannot, to our knowledge, be estimated using standard generalized linear mixed models. For these reasons, any expert on biometrical genetics should be familiar with both approaches and choose the one that is most suitable for a given problem.

It should be noted that likelihood ratio tests cannot be based on the conventional $\chi^2$ distribution when testing variance components. In the context of twin models Dominicus

et al. (2005) derive the correct null distributions that are mixtures of $\chi^2$ distributions.

We have illustrated how statistical equivalence between models can be exploited to reduce the number of random effects and hence computational complexity. Rabe-Hesketh and Skrondal (2001) and Skrondal and Rabe-Hesketh (2004) discuss other examples where the number of random effects can be reduced by reparameterization.

The models discussed in this article are straightforward to extend in several ways, still using standard software for mixed models. For the family designs discussed in Section 5 we could also include dominance and maternal effects as nicely demonstrated by Pawitan et al. (2004). In addition to the kinds of covariates considered in Section 7.1, we could include covariates representing measured genotypes at specific loci (e.g. Burton et al., 1999; van den Oord, 2001). These effects can be made random at one or more levels to study interactions between genotype and environmental components ('genotype–environment interaction') or interactions between genotype and genetic components ('epistasis'). Higher-level random effects could be used to model unobserved environmental heterogeneity at different levels, for instance between neighborhoods.

## 9. Supplementary Materials

Web-based supplementary materials, including example data sets and corresponding commands for `Stata` and `R/S-PLUS`, are available under the Paper Information link at the *Biometrics* website `http://www.tibs.org/biometrics`.

### REFERENCES

Breslow, N. E. (2005). Whither PQL? In *Proceedings of the Second Seattle Symposium in Biostatistics: Analysis of Correlated Data*, D. Y. Lin and P. J. Heagerty (eds), 1–22. New York: Springer.

Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association* **88,** 9–25.

Browne, W. J., Rasbash, J., and Ng, E. W. S. (2005). MCMC estimation in MLwiN version 2.0. Bristol, U.K.: University of Bristol. Downloadable from `http://www.cmm.bristol.ac.uk/MLwiN/download/manuals.shtml`.

Burton, P. R., Tiller, K. J., Currin, L. C., Cookson, W. O. C. M., Musk, A. W., and Palmer, L. J. (1999). Genetic variance components analysis for binary phenotypes using generalized linear mixed models (GLMMs) and Gibbs sampling. *Genetic Epidemiology* **17,** 118–140.

Dominicus, A., Skrondal, A., Gjessing, H. K., Pedersen, N., and Palmgren, J. (2005). Likelihood ratio tests in behavioral genetics: Problems and solutions. *Behavior Genetics* **36,** 331–340.

Elston, R. C. (2001). Introduction and overview. *Statistical Methods in Medical Research* **9,** 527–541.

Falconer, D. S. and MacKay, T. F. C. (1996). *Introduction to Quantitative Genetics*, 4th edition. Essex: Longman.

Fisher, R. A. (1918). The correlation between relatives on the supposition of Mendelian inheritance. *Transactions of the Royal Society of Edinburgh* **52,** 399–433.

Guo, G. and Wang, J. (2002). The mixed or multilevel model for behavior genetic analysis. *Behavior Genetics* **32,** 37–49.

Hedeker, D. and Gibbons, R. D. (1996). MIXOR: A computer program for mixed-effects ordinal probit and logistic regression analysis. *Computer Methods and Programs in Biomedicine* **49,** 157–176.

Jinks, J. L. and Fulker, D. W. (1970). Comparison of the biometrical genetical, MAVA, and classical approaches to the analysis of human behavior. *Psychological Bulletin* **73,** 311–349.

Khoury, M. J., Beaty, T. H., and Cohen, B. C. (1993). *Fundamentals of Genetic Epidemiology.* Oxford: Oxford University Press.

Lillard, L. A. and Panis, C. W. A. (2000). *aML User's Guide and Reference Manual.* Los Angeles, California: EconWare.

Magnus, P., Gjessing, H. K., Skrondal, A., and Skjærven, R. (2001). Paternal contribution to birth weight. *Journal of Epidemiology and Community Health* **55,** 873–877.

McArdle, J. J., and Prescott, C. A. (2005). Mixed-effects variance components models for biometric family analyses. *Behavior Genetics* **35,** 631–652.

Neale, M. C. and Maes, H. H. (2004). *Methodology for Genetic Studies of Twins and Families.* Richmond, VA: Virginia Commonwealth University, Department of Psychiatry. Downloadable from `http://ibgwww.colorado.edu/workshop2004/cdrom/HTML/book2004a.pdf`.

Neale, M. C., Boker, S. M., Xie, G., and Maes, H. H. (2004). *Mx: Statistical Modeling (Sixth Edition).* Richmond, Virginia: Virginia Commonwealth University, Department of Psychiatry. Downloadable from `http://www.vipbg.vcu.edu/mxgui/`.

Pawitan, Y., Reilly, M., Nilsson, E., Cnattingius, S., and Lichtenstein, P. (2004). Estimation of genetic and environmental factors for binary traits using family data. *Statistics in Medicine* **23,** 449–465.

Rabe-Hesketh, S. and Skrondal, A. (2001). Parameterization of multivariate random effects models for categorical data. *Biometrics* **57,** 1256–1264.

Rabe-Hesketh, S. and Skrondal, A. (2005). *Multilevel and Longitudinal Modeling Using Stata.* College Station, Texas: Stata Press.

Rabe-Hesketh, S. and Skrondal, A. (2007). Generalized linear mixed effects models. In *Longitudinal Data Analysis: A Handbook of Modern Statistical Methods*, G. Fitzmaurice, M. Davidian, G. Mdenberghs, and G. Verbeke (eds). Boca Raton, Florida: Chapman & Hall/CRC.

Rabe-Hesketh, S., Skrondal, A., and Pickles, A. (2005). Maximum likelihood estimation of limited and discrete dependent variable models with nested random effects. *Journal of Econometrics* **128,** 301–323.

Rodriguez, G. and Goldman, N. (1995). An assessment of estimation procedures for multilevel models with binary responses. *Journal of the Royal Statistical Society, Series A* **158,** 73–89.

Rodriguez, G. and Goldman, N. (2001). Improved estimation procedures for multilevel models with binary response: A case study. *Journal of the Royal Statistical Society, Series A* **164,** 339–355.

Skrondal, A. and Rabe-Hesketh, S. (2004). *Generalized Latent Variable Modeling: Multilevel, Longitudinal, and Structural Equation Models.* Boca Raton, Florida: Chapman & Hall/CRC.

Spiegelhalter, D. J., Thomas, A., Best, N. G., and Gilks, W. R. (1996). *BUGS 0.5 Bayesian Analysis using Gibbs Sampling. Manual (version ii).* Cambridge: MRC-Biostatistics Unit. Downloadable from `http://www.mrc-bsu.cam.ac.uk/bugs/documentation/contents.shtml`.

van den Oord, E. J. C. G. (2001). Estimating effects of latent and measured genotypes in multilevel models. *Statistical Methods in Medical Research* **10,** 393–407.